

Contents

1	Q1:Hadoop Ecosystem Technologies Discussion (10 points)	2
2	Q2:Pig Script Commands and Data Processing (10 points)	3
2.1	Pig Code	3
2.1.1	Q2 SnapShot	4
3	Q3:Hive Queries (10 points)	5
3.1	HIVE Code	5
3.1.1	Q3 SnapShot	6
4	Q4:Weblog Data Analysis (10 points)	7
4.1	Q4 Code	7
4.1.1	Q4 SnapShot	7

1 Q1:Hadoop Ecosystem Technologies Discussion (10 points)

1. **YARN (Yet Another Resource Negotiator):** YARN is crucial for resource management in Hadoop. It allocates system resources (like CPU and memory) to various running applications. It improves cluster efficiency and allows for simultaneous data processing, making it vital for managing big data workloads.
2. **Zookeeper:** This is a coordination service for distributed applications. In the context of Hadoop, Zookeeper maintains configuration information, provides distributed synchronization, and manages a cluster's naming registry. It ensures high availability and reliability of the system, crucial for big data applications that can't afford downtime or data inconsistency.
3. **Oozie:** Oozie is Hadoop's scheduler. It's used to manage and schedule Hadoop jobs, such as MapReduce and Pig. Oozie simplifies complex job workflows, allowing you to define a series of jobs to be executed in a specific order. It's essential for automating and optimizing big data processes, ensuring tasks are performed efficiently and on schedule.
4. **Sqoop/Hue:** Sqoop is a tool designed to transfer data between Hadoop and relational databases. It's crucial for integrating big data with traditional data warehouses, allowing for efficient data import/export. Hue, on the other hand, is a web interface for Hadoop services. It simplifies working with Hadoop, providing a user-friendly interface for interacting with various components like HDFS, MapReduce, and YARN. Hue makes it easier to manage big data workflows, especially for those who might not be as comfortable with command-line operations.

2 Q2:Pig Script Commands and Data Processing (10 points)

2.1 Pig Code

```
literate
-- Load daily stock data
daily = LOAD "s3://bigdata-hw6-lanston/pig/NYSE_daily" USING
    PigStorage(',') AS AS (
        exchange: chararray,
        stock: chararray,
        date: chararray,
        open_price: float,
        high_price: float,
        low_price: float,
        close_price: float,
        volume: long,
        adj_close: float
    );

-- Load dividends data
dividends = LOAD "s3://bigdata-hw6-lanston/pig/NYSE_dividends"
    USING PigStorage(',') AS (
        exchange: chararray,
        stock: chararray,
        date: chararray,
        dividend: float
    );

-- Join the datasets on stock and date
joinedData = JOIN daily BY (stock, date), dividends BY (stock,
    date);

-- Calculate dividend/close_price
calculatedData = FOREACH joinedData GENERATE daily::stock AS
    stock, daily::date AS date,
    dividends::dividend/daily::close_price AS div_close_ratio;

-- Group all
calculatedData2 = GROUP calculatedData ALL;

-- Calculate the min and max ratios
calculatedData3 = FOREACH calculatedData2 GENERATE
    MIN(dividend_ratio.ratio) AS min_ratio,
```

```

MAX(dividend_ratio.ratio) AS max_ratio;

-- Join to find the records with min and max ratios
min_record = JOIN joinedData BY ratio, min_max_ratios BY
    min_ratio;
max_record = JOIN joinedData BY ratio, min_max_ratios BY
    max_ratio;

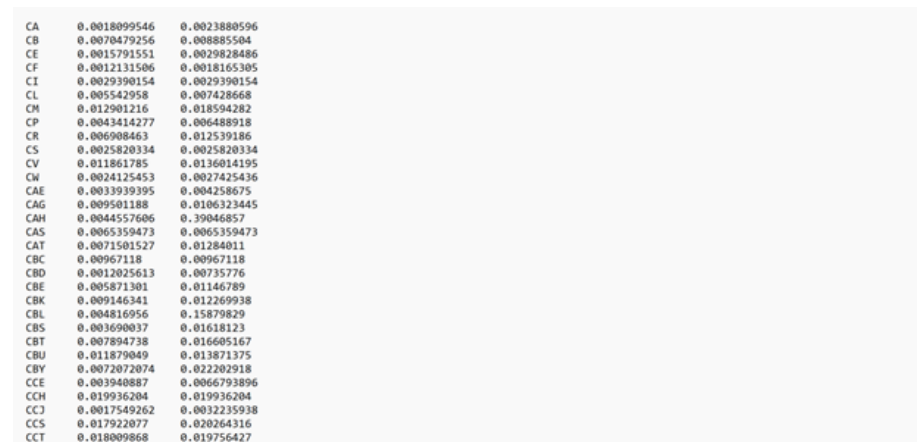
-- Prepare the final records for min and max
min_record_final = FOREACH min_record GENERATE
    FLATTEN(dividend_ratio::stock) AS stock,
    FLATTEN(dividend_ratio::date) AS date,
    dividend_ratio::ratio;

max_record_final = FOREACH max_record GENERATE
    FLATTEN(dividend_ratio::stock) AS stock,
    FLATTEN(dividend_ratio::date) AS date,
    dividend_ratio::ratio;

-- Store the final results
STORE min_record_final INTO
    's3://bigdata-hw6-lanston/pig/pig_min' USING PigStorage(',');
STORE max_record_final INTO
    's3://bigdata-hw6-lanston/pig/pig_max' USING PigStorage(',');

```

2.1.1 Q2 SnapShot



CA	0.0018099546	0.0023880596
CB	0.0070479256	0.008885504
CE	0.0015791551	0.0029828486
CF	0.0012131506	0.0018165305
CI	0.0029390154	0.0029390154
CL	0.005542958	0.007428668
CM	0.012901216	0.018594282
CP	0.0043414277	0.006488918
CR	0.006908463	0.012539186
CS	0.0025820334	0.0025820334
CV	0.011861785	0.0136014195
CW	0.0024125453	0.0027425436
CAE	0.0033939395	0.004258675
CAG	0.009501188	0.0106323445
CAH	0.0044557606	0.39046857
CAS	0.0065359473	0.0065359473
CAT	0.0071501527	0.01284011
CBC	0.00967118	0.00967118
CBD	0.0012025613	0.00735776
CBE	0.005871301	0.01146789
CBK	0.009146341	0.012269938
CBL	0.004816956	0.15879829
CBS	0.003690037	0.01618123
CBT	0.007894738	0.016605167
CBU	0.011879049	0.013871375
CBY	0.0072072074	0.022202918
CCE	0.003940887	0.0066793896
CCM	0.019936204	0.019936204
CCJ	0.0017549262	0.0032235938
CCS	0.017922077	0.020264316
CCT	0.018009868	0.019756427

Figure 1: [pig shell

3 Q3:Hive Queries (10 points)

3.1 HIVE Code

```
literate
-- Create an external table nyTaxi
CREATE EXTERNAL TABLE nyTaxi (
    VendorID INT,
    lpep_pickup_datetime STRING,
    lpep_dropoff_datetime STRING,
    store_and_fwd_flag STRING,
    RatecodeID INT,
    PULocationID INT,
    DOLocationID INT,
    passenger_count INT,
    trip_distance FLOAT,
    fare_amount FLOAT,
    extra FLOAT,
    mta_tax FLOAT,
    tip_amount FLOAT,
    tolls_amount FLOAT,
    ehail_fee STRING,
    improvement_surcharge FLOAT,
    total_amount FLOAT,
    payment_type INT,
    trip_type INT
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
LINES TERMINATED BY '\n'
STORED AS TEXTFILE
LOCATION 's3://bigdata-hw6-lanston/hive/';

-- Get distinct RatecodeID from the table
SELECT DISTINCT RatecodeID FROM nyTaxi;

-- Show all rows/columns where RatecodeID = 1
SELECT * FROM nyTaxi WHERE RatecodeID = 1;
```

3.1.1 Q3 SnapShot

Penology V12.31.40.064															
U/V1/3/5023	U/V1/3/5023	W	1	130	102	1	2,440	9,50	9,50	9,50	0,00	0,00	0,00	10,0	2
U/V1/3/5046	U/V1/3/5028	W	1	130	100	1	21,291	40,50	9,50	9,50	0,00	0,00	0,00	41,8	2
U/V1/3/5057	U/V1/3/5028	W	1	130	100	1	2,440	9,50	9,50	9,50	0,00	0,00	0,00	10,0	2
U/V1/3/5045	U/V1/3/5057	W	1	130	205	1	7,001	21,00	9,50	9,50	0,00	0,00	0,00	22,5	1
U/V1/3/5054	U/V1/3/5054	W	1	130	102	1	1,102	4,00	9,50	9,50	0,00	0,00	0,00	5,0	2
U/V1/3/5020	U/V1/3/5020	W	1	130	115	1	1,102	11,00	9,50	9,50	0,00	0,00	0,00	12,5	2
U/V1/3/5081	U/V1/3/5057	W	1	36	190	1	1,102	6,00	9,50	9,50	0,00	0,00	0,00	7,0	2
U/V1/3/5020	U/V1/3/5057	W	1	36	175	1	1,102	5,50	9,50	9,50	0,00	0,00	0,00	6,50	2
U/V1/3/5038	U/V1/3/5046	W	1	17	37	1	1,109	9,00	9,50	9,50	0,00	0,00	0,00	9,00	2
U/V1/3/5046	U/V1/3/5046	W	1	230	109	1	4,700	17,50	9,50	9,50	0,00	0,00	0,00	18,00	2
U/V1/3/5044	U/V1/3/5044	W	1	230	109	1	1,200	7,00	9,50	9,50	1,74	9,00	0,00	13,50	1
U/V1/3/5057	U/V1/3/5011	W	1	41	41	1	0,351	4,50	9,50	9,50	0,00	0,00	0,00	5,0	2
U/V1/3/5032	U/V1/3/5011	W	1	41	41	1	89,5	9,50	9,50	9,50	0,00	0,00	0,00	90,0	2
U/V1/3/5081	U/V1/3/5011	W	1	229	41	1	2,066	9,50	9,50	9,50	2,7	9,00	0,00	13,5	1
U/V1/3/5057	U/V1/3/5011	W	1	42	42	1	0,351	4,50	9,50	9,50	0,00	0,00	0,00	5,0	2
U/V1/3/5021	U/V1/3/5034	W	1	101	14	2	5,300	17,50	9,50	9,50	1,44	9,00	0,00	14,44	1
U/V1/3/5035	U/V1/3/5034	W	1	230	223	1	15,49	47,50	9,50	9,50	0,00	0,00	0,00	47,0	2
U/V1/3/5023	U/V1/3/5034	W	1	129	129	1	5,300	17,50	9,50	9,50	1,44	9,00	0,00	14,44	1
U/V1/3/5046	U/V1/3/5051	W	1	230	240	0	8,00	6,50	9,50	9,50	0,00	0,00	0,00	7,0	2
U/V1/3/5052	U/V1/3/5051	W	1	49	240	1	2,100	14,00	9,50	9,50	0,00	0,00	0,00	14,0	2
U/V1/3/5031	U/V1/3/5042	W	1	49	94	1	2,300	15,50	9,50	9,50	1,34	9,00	0,00	15,34	1
U/V1/3/5047	U/V1/3/5051	W	1	149	47	1	1,17	7,50	9,50	9,50	0,00	0,00	0,00	8,0	2
U/V1/3/5034	U/V1/3/5034	W	1	133	9	1	5,271	17,50	9,50	9,50	1,44	9,00	0,00	14,44	1
U/V1/3/5042	U/V1/3/5059	W	1	255	40	1	8,00	27,50	9,50	9,50	0,00	0,00	0,00	28,0	2
U/V1/3/5059	U/V1/3/5059	W	1	255	40	1	14	6,00	9,50	9,50	0,00	0,00	0,00	6,00	2
U/V1/3/5039	U/V1/3/5018	W	1	114	110	1	9,20	27,00	9,50	9,50	0,00	0,54	0,00	33,54	2
U/V1/3/5059	U/V1/3/5018	W	1	230	259	1	14	6,00	9,50	9,50	0,00	0,00	0,00	6,00	2
U/V1/3/5059	U/V1/3/5018	W	1	92	56	1	1,40	6,00	9,50	9,50	0,00	0,00	0,00	6,00	2
U/V1/3/5036	U/V1/3/5046	W	1	129	240	1	1,4	7,00	9,50	9,50	0,00	0,00	0,00	7,0	2
U/V1/3/5036	U/V1/3/5046	W	1	130	215	1	2,15	11,00	9,50	9,50	0,00	0,00	0,00	11,00	2
U/V1/3/5058	U/V1/3/5046	W	1	255	226	1	2,7	11,50	9,50	9,50	4,00	9,00	0,00	16,5	1
U/V1/3/5058	U/V1/3/5046	W	1	230	122	1	12,2	39,50	9,50	9,50	0,00	0,00	0,00	39,5	2
U/V1/3/5059	U/V1/3/5025	W	1	42	231	1	12,00	38,50	9,50	9,50	0,00	0,00	0,00	38,5	2
U/V1/3/5056	U/V1/3/5058	W	1	129	251	1	5,00	14,50	9,50	9,50	2,00	9,00	0,00	19,5	1
U/V1/3/5056	U/V1/3/5058	W	1	121	96	1	1,10	6,00	9,50	9,50	0,00	0,00	0,00	6,00	2
U/V1/3/5056	U/V1/3/5058	W	1	114	243	1	2,0	11,00	9,50	9,50	0,00	0,00	0,00	11,0	2
U/V1/3/5056	U/V1/3/5058	W	1	129	196	1	1,10	6,00	9,50	9,50	0,00	0,00	0,00	6,00	2
U/V1/3/5056	U/V1/3/5052	W	1	41	238	1	1,4	7,00	9,50	9,50	1,10	9,00	0,00	9,10	1
U/V1/3/5056	U/V1/3/5056	W	1	14	230	1	1,10	15,00	9,50	9,50	0,00	0,00	0,00	16,0	2
U/V1/3/5056	U/V1/3/5052	W	1	41	40	1	0,35	4,50	9,50	9,50	0,00	0,00	0,00	4,5	2
U/V1/3/5056	U/V1/3/5056	W	1	41	42	1	2,1	9,00	9,50	9,50	0,00	0,00	0,00	10,0	1
U/V1/3/5056	U/V1/3/5014	W	1	41	10	1	0,3	1,00	9,50	9,50	0,00	0,00	0,00	1,0	2
U/V1/3/5058	U/V1/3/5028	W	1	41	90	1	13,0	40,00	9,50	9,50	0,00	0,00	0,00	41,00	2
U/V1/3/5058	U/V1/3/5028	W	1	90	40	1	1,10	6,00	9,50	9,50	0,00	0,00	0,00	6,00	2
U/V1/3/5058	U/V1/3/5028	W	1	240	252	1	6,7	25,50	9,50	9,50	0,00	0,00	0,00	27,0	2
U/V1/3/5034	U/V1/3/5046	W	1	80	36	1	2,4	10,00	9,50	9,50	2,25	9,00	0,00	13,50	1
U/V1/3/5054	U/V1/3/5054	W	1	1	23	1	2,7	11,5	9,50	9,50	0,00	0,00	0,00	11,5	2
U/V1/3/5058	U/V1/3/5054	W	1	108	237	1	2,2	9,5	9,50	9,50	1,10	9,00	0,00	12,90	1
U/V1/3/5057	U/V1/3/5054	W	1	130	102	1	2,440	9,50	9,50	9,50	0,00	0,00	0,00	10,0	2
U/V1/3/5033	U/V1/3/5015	W	1	133	70	1	9,5	26,00	9,50	9,50	0,00	1,34	0,00	31,34	2
U/V1/3/5052	U/V1/3/5057	W	1	104	151	1	1,3	6,00	9,50	9,50	0,00	0,00	0,00	6,0	2
U/V1/3/5052	U/V1/3/5057	W	1	137	36	1	0,7	2,50	9,50	9,50	0,00	0,00	0,00	2,5	2
U/V1/3/5052	U/V1/3/5036	W	1	230	223	2	21,7	47,00	9,50	9,50	18,50	0,00	0,00	49,5	1
U/V1/3/5052	U/V1/3/5036	W	1	17	2	1	0,02	0,50	9,50	9,50	0,00	0,00	0,00	0,5	2
U/V1/3/5052	U/V1/3/5042	W	1	255	17	1	25	18,00	9,50	9,50	2,50	9,00	0,00	22,0	1
U/V1/3/5052	U/V1/3/5046	W	1	104	80	1	6,40	20,00	9,50	9,50	5,30	9,00	0,00	26,40	2
U/V1/3/5030	U/V1/3/5030	W	1	40	43	1	12,5	37,50	9,50	9,50	0,00	0,00	0,00	37,5	2

(Line Labels: 2,15 seconds, Pathology 15000 rows)

Figure 2: [hive shell

4 Q4:Weblog Data Analysis (10 points)

4.1 Q4 Code

```
literate
CREATE EXTERNAL TABLE tripadvisor_logs (
    'ip' STRING,
    'timestamp' STRING,
    'request' STRING,
    'status' INT,
    'bytes' BIGINT,
    'referrer' STRING,
    'useragent' STRING
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
LINES TERMINATED BY '\n'
STORED AS TEXTFILE
LOCATION 'hdfs:///hive/bigdata/hw6/tripadvisor';

SELECT ip, COUNT(*) AS error_count
FROM tripadvisor_logs
WHERE status = 404
GROUP BY ip;
```

4.1.1 Q4 SnapShot

```
> GROUP BY ip;
Query ID = hadoop_20231115003150_e6a89465-ddcf-4bab-a0ef-333353789d2d
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1700004796704_0005)

-----
VERTICES      MODE           STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
-----
Map 1 ..... container  SUCCEEDED    1         1         0         0         0         0
Reducer 2 ..... container  SUCCEEDED    2         2         0         0         0         0
-----
VERTICES: 02/02 [=====>>] 100% ELAPSED TIME: 8.20 s
```

Figure 3: [shell